

# Identifiability of a Coalescent-based Population Tree Model

Arindam RoyChoudhury

April 15, 2013

## Abstract

Identifiability of evolutionary tree models has been a recent topic of discussion and some models have been shown to be non-identifiable. A coalescent-based rooted population tree model, originally proposed by Nielsen et al. 1998 [2], has been used by many authors in the last few years and is a simple tool to accurately model the changes in allele frequencies in the tree. However, the identifiability of this model has never been proven. Here we prove this model to be identifiable by showing that the model parameters can be expressed as functions of the probability distributions of subsamples. This is a step toward proving the consistency of the maximum likelihood estimator of the population tree based on this model.

## 1 Introduction

A rooted evolutionary tree is a directed weighted tree graph; it represents the evolutionary relationship between groups (also called taxa) of organisms (Figure 1(a)). A leaf or a tip is a node with degree 1; each tip represents a modern day taxon. The root (node 0) represents the most recent common ancestor (MRCA) of all the taxa. The direction (of evolution) is from the root to the tips. Evolutionary tree as a vector of parameters influences the probability distribution of alleles at the tips.

A rooted population tree is a rooted evolutionary tree where the taxa are populations from the same species. Two types of parameters are common in any model of the rooted population tree: the tree-topology parameter (a categorical parameter) for the whole tree, and a branch parameter for each branch (also called edge).

The tree-topology is the order in which the path from the root separates for the given set of populations; it is represented as a directed tree graph without the weight. (In Figure 1(a) and (b), the two trees have different tree-topologies for the populations 1-4.) A branch parameter is usually a branch-length (an edge-weight) or a transition probability matrix that influences the change in allele frequency between the two nodes of a branch.

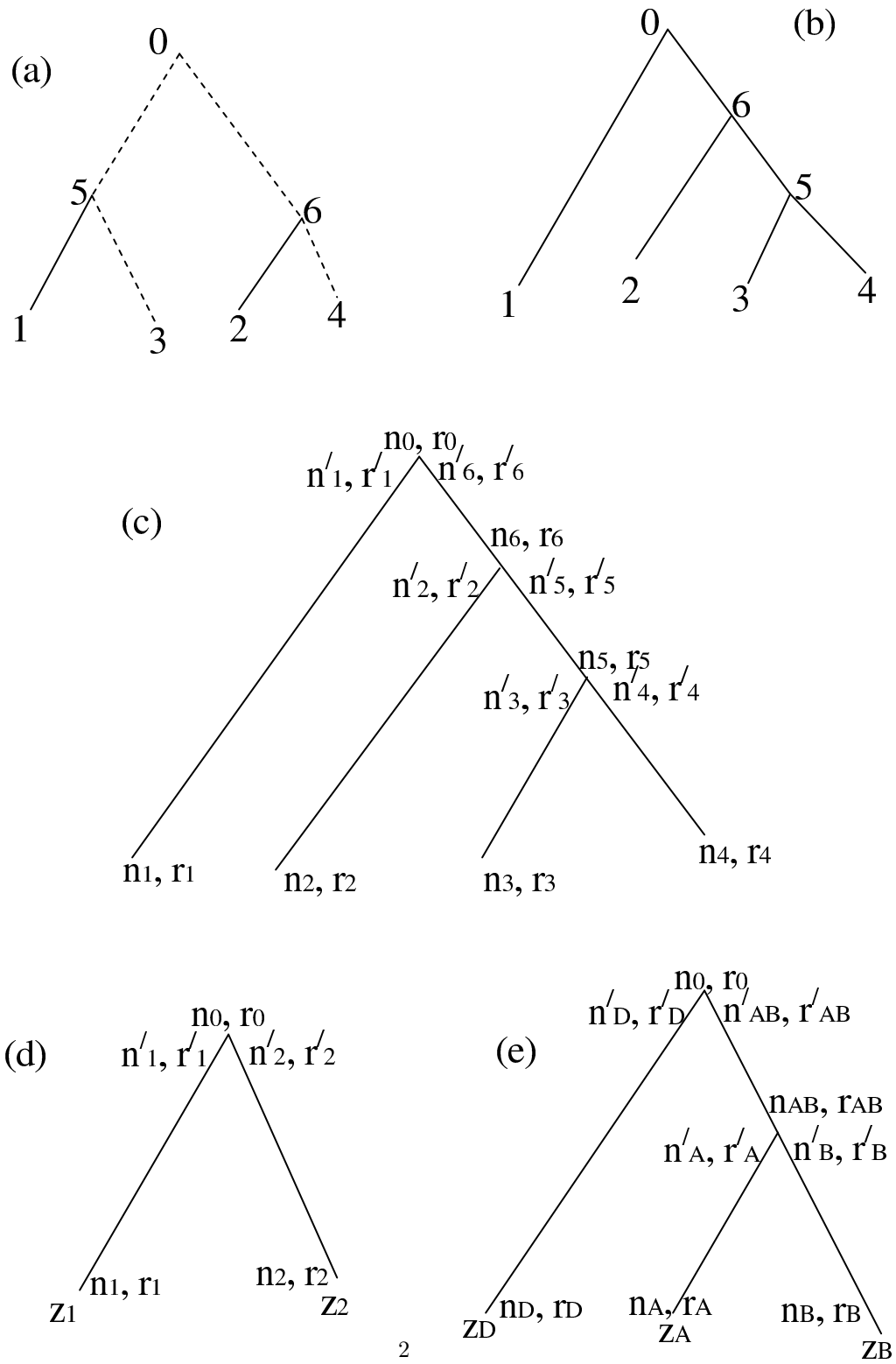


Figure 1: Population trees

Here we will prove the identifiability of a population tree model by [2, 5] that uses Kingman’s Coalescent Process ([3]). The model was later modified and expanded by various authors ([4, 5, 6, 7, 8]). Coalescent-based models are of significant importance as they model the underlying allele frequency changes with accuracy and relative ease (see [13]).

Due to the underlying structure in evolutionary tree-based models, its identifiability is never obvious. The identifiability of certain evolutionary tree models have been a recent topic of discussion. [1] proved the identifiability of a general time reversible (GTR) transition probability matrix-based model. Non-identifiability of another time reversible model was established in [10]. The non-identifiability of mixture models have been discussed in [11]. The identifiability for the [12] model has been proven by [9]. To our knowledge the identifiability of the coalescent-based model of [2, 5] has never been proven.

For estimating evolutionary trees each independent genetic locus is viewed as a single data-point, as opposed to viewing each individual as a data-point (see, for example [6]). Thus, identifiability would mean that the model parameters can be identified from the distribution of allele-types for a set of individuals at a single genetic locus.

## 2 The model

In this section we will describe the underlying model of [2, 5]. We start by defining our notations (see also Figure 1(c)). We define a  $P$ -tip population tree as  $T = (\Lambda^{(P)}, \Psi, \theta)$ . The parameter  $\Lambda^{(P)}$  is the tree-topology, an unweighted directed tree-graph; it takes finitely many discrete categorical values; the  $(P)$  in superscript denotes the number of tips. The parameter  $\Psi = (\tau_1, \tau_2, \dots, \tau_{2P-2})$  is a vector of length  $2P - 2$  consisting of the branch-lengths  $\tau_i$  for each branch  $i$  in  $\Lambda^{(P)}$ . A strictly bifurcating tree-topology has exactly  $2P - 2$  branches. If  $\Lambda^{(P)}$  is non-bifurcating then it has less branches and the remaining elements of  $\Psi$  are populated by zeros. The parameter  $\theta$  is a vector containing the parameters of root distribution which we will define later in this section. We also define  $S(\Lambda^{(P)})$  as the set of tips at  $\Lambda^{(P)}$ .

At each tip  $z$  there are  $n_z (\geq 1)$  lineages, each having allele-type ‘0’ or ‘1’. The allele types among these lineages at each tip are the observable random variables. Similarly, at each non-tip node  $x$ , the random variable  $n_x (\geq 1)$  is the (random) number of lineages that are ancestral to the tips below  $x$  along the tree. We also define the random variable  $r_x$  at each node  $x$  (tip or non-tip), as the count of allele ‘1’ among the  $n_x$  lineages. From now on we will use the term ‘allele-count’ to refer to the count of allele ‘1’. For each tip  $z$ , the allele-count  $r_z$  is observable.

Consider a branch with lower (towards the tips) node  $x$  and upper (towards the root) node  $y$ . Let  $n'_x$  be the number of lineages in  $y$  that are ancestral to the  $n_x$  lineages at  $x$  ( $n'_x \leq n_x$ ). Also, let  $r'_x$  be the allele-count among these  $n'_x$  lineages ( $r'_x \leq r_x$ ). If  $y$  is the upper node of  $\nu$  branches with lower nodes

$x_1, x_2, \dots, x_\nu$ , then

$$n_{x_1}, n_{x_2}, \dots, n_{x_\nu} \text{ are independent, and } n_y = \sum_{k=1}^{\nu} n'_{x_k} \quad (1)$$

and also  $r_y = \sum_{k=1}^{\nu} r'_{x_k}$ . (For a strictly bifurcating tree  $\nu = 2$ .)

From the model parameters  $T = (\Lambda^{(P)}, \Psi, \theta)$  one computes the probability of observed vector of allele-counts  $\mathbf{r} = (r_1, r_2, \dots, r_P)$  from samples of sizes  $\mathbf{n} = (n_1, n_2, \dots, n_P)$  at  $P$  tips  $(1, 2, \dots, P)$  as follows. Consider a branch with length  $\tau_{x_1}$ , with upper node  $y$  and lower node  $x_1$ . Given the probability mass function (pmf) of  $n_{x_1}$  (the number of lineages at  $x_1$ ), the pmf of  $n'_{x_1}$  is computed as

$$\Pr_{\mathbf{n}}(n'_{x_1} = i' | n_{x_1} = i; \tau_{x_1}) = \left( \prod_{j=i'+1}^i \lambda_j \right) \sum_{j=i'}^i \frac{e^{-\lambda_j \tau_{x_1}}}{\prod_{j'=i', j' \neq j}^i (\lambda_{j'} - \lambda_j)}, \quad (2)$$

where  $\lambda_j = j(j-1)/2$ . Then, the pmf of  $n_y$  is determined from Eq. (1).

Using Eqs. (2) and (1), starting from  $\mathbf{n} = (n_1, n_2, \dots, n_P)$  and going upward, one computes the pmf of  $n_z$  and  $n'_z$  for any non-tip non-root node  $z$ , and finally  $n_0$  at the root (node 0). Then a ‘root distribution’ with parameter  $\theta$  gives the pmf of (allele-count)  $r_0$  given  $n_0$  at the root:

$$G(0) = \left( \Pr_{\mathbf{n}}(r_0 = j | n_0 = i; \theta), j = 0, 1, \dots, n_0; i = 1, 2, \dots, m_0^{(b)} \right),$$

where

$$m_0^{(b)} = \sum_{z \text{ is a tip}} n_z$$

is the maximum possible value of  $n_0$  (number of lineages at the root). Different authors have used different root distributions. In particular [5] used symmetric Beta-Binomial distribution:

$$\Pr_{\mathbf{n}}(r_0 = j | n_0 = i; \theta) = \binom{i}{j} \frac{\beta(j + \theta) \beta(i - j + \theta)}{\beta(\theta, \theta)}, \quad (3)$$

where  $\beta(.,.)$  is the Beta Function;  $\theta > 0$  is a parameter to be estimated.

Then, from the distribution of  $n_0, r_0$  and  $(n_z, n'_z)$  for all non-root nodes  $z$ , we compute the distribution of  $r_z$  (allele-counts) at the rest of the nodes as follows. Consider a node  $y$  where  $\nu$  branches merge from the bottom with the bottom nodes  $x_1, x_2, \dots, x_\nu$ . Recall that we already have the distributions of  $n_y, n_{x_i}$  and  $n'_{x_i}$ ,  $i = 1, 2, \dots, \nu$ . The pmf of  $r'_{x_i}$  is computed from the pmf of  $r_y$  using the formula

$$\begin{aligned} & \Pr_{\mathbf{n}}(r'_{x_1} = j'_1, r'_{x_2} = j'_2, \dots, r'_{x_\nu} = j'_\nu | r_y = j, n_y = i, n'_{x_1} = i'_1, n'_{x_2} = i'_2, \dots, n'_{x_\nu} = i'_\nu) \\ &= \frac{\binom{j}{j'_1, j'_2, \dots, j'_\nu} \binom{i-j}{i'_1 - j'_1, i'_2 - j'_2, \dots, i'_\nu - j'_\nu}}{\binom{i}{i'_1, i'_2, \dots, i'_\nu}}. \end{aligned} \quad (4)$$

Then the pmf of  $r_{x_k}$  is computed from the above pmf using the following (from an expression in [5]):

$$\begin{aligned} & \Pr_{\mathbf{n}}(r_{x_k} = j_k \mid r'_{x_k} = j'_k, n'_{x_k} = i'_k, n_{x_k} = i_k) \\ &= \frac{\beta(j_k, i_k - j_k)}{\beta(j'_k, i'_k - j'_k)} \binom{i_k - i'_k}{j_k - j'_k}, 0 < j_k < i_k \text{ and } 0 < j'_k < i'_k, \\ & 1, 0 = j_k = j'_k \text{ or } 0 = i_k - j_k = i'_k - j'_k, \\ & 0, \text{ otherwise;} \end{aligned} \quad (5)$$

$k = 1, 2, \dots, \nu$  ([5]). Thus, starting with  $G(0)$  at the root, one computes the joint pmf of  $(r_1, r_2, \dots, r_P)$  from the formulae in Eqs. (4) and (5). Note that in Eqs. (1), (2), (3), (4) and (5) probability ‘flows’ up along  $n$ ’s and then flows down along  $r$ ’s.

Now that we have completely described the model, we will proceed to prove the identifiability of this model in the next section.

### 3 Identifiability

Let  $T = (\Lambda^{(P)}, \Psi, \theta)$  be a tree with  $S(\Lambda^{(P)}) = \{1, 2, \dots, P\}$ . We define a subtree  $T^*$  of  $T$  as a tree formed by a subset  $S^*$  (cardinality  $P' \leq P$ ) of  $S(\Lambda^{(P)})$  by tracking the tips in  $S^*$  along the tree to their most recent common ancestor (MRCA) node. Thus,  $T^* = (\Lambda^{(P')*}, \Psi^*, \theta)$ , where  $\Lambda^{(P')*}$  is the tree-topology with  $P'$  tips of  $S^*$ . For example, in Figure 1(a),  $P = 4$ ,  $S(l4) = \{1, 2, 3, 4\}$ ,  $S^* = \{3, 4\}$  and the subtree  $T^*$  is drawn with the dotted lines.

Consider two distinct trees  $T_1 = (\Lambda_1^{(P)}, \Psi_1, \theta_1)$  and  $T_2 = (\Lambda_2^{(P)}, \Psi_2, \theta_2)$  with a common set of tips  $S_{T_1,2} = S(\Lambda_1^{(P)}) = S(\Lambda_2^{(P)})$ .

If  $\theta_1 = \theta_2 = \theta$ , then there must be at least one doubleton subset  $\{z_1, z_2\} \subseteq S_{T_1,2}$  with the following property: the subtrees  $T_1^* = (\Lambda^{(2)}, \Psi_1^*, \theta)$  and  $T_2^* = (\Lambda^{(2)}, \Psi_2^*, \theta)$ , formed by tracking  $z_1$  and  $z_2$  to the root in  $T_1$  and  $T_2$  (respectively), are distinct. That is, if  $\Psi_l^* = (\tau_{1l}, \tau_{2l})$  and  $\tau_{jl}$  is the path distance (total branch length) between  $z_j$  and the MRCA of  $z_1$  and  $z_2$  along the subtree  $T_l^*$  ( $j, l = 1, 2$ ), then  $(\tau_{11}, \tau_{21}) \neq (\tau_{12}, \tau_{22})$ . (Note that there is only one possible tree-topology for a two-tip tree, denoted as  $\Lambda^{(2)}$  above.) Thus, the set of all two tip subtrees, along with  $\theta$ , uniquely identifies the tree.

We assign the two-tip subtrees into two categories: Type-I subtrees are those with the root as the MRCA of the two tips. For example in Figure 1(a), the subtree formed by tips  $\{3, 4\}$  has the root as the MRCA of the two tips 3 and 4. Thus, it is of Type-I. All other two-tip subtrees are Type-II subtrees. For example, in Figure 1(a), if a subtree is formed by tips 2 and 4, it will be a Type-II subtree as their MRCA is node 6, and not the root. We will deal with these two types of subtrees separately.

We note that the root distribution of [5] (Eq. (3)) is identifiable as it is Beta-Binomial. Next, we will prove the identifiability of the whole model by assuming a general identifiable root distribution that has parameter vector  $\theta$ . (In particular, our proof would work with Beta-Binomial as the root distribution.)

**Theorem** Suppose that we have a tree  $T$  with the underlying model as described in Section 2. Also, suppose that we have  $N_k \geq 2$  lineages sampled at each tip  $k$  and the root distribution is identifiable. Then the parameters of  $T$  are identifiable from the distribution of allele types at the tips.

To prove the above theorem, we will show that the parameters of each two-tip subtree can be expressed as a function of the joint pmf

$$(\Pr_{\mathbf{n}}((R_1, R_2, \dots, R_P) = (J_1, J_2, \dots, J_P); T), J_k = 0, 1, 2, \dots, N_k, k = 1, 2, \dots, P). \quad (6)$$

This will complete the proof as the set of all two-tip subtrees, along with  $\theta$ , uniquely identifies the tree.

### 3.1 Identifiability of Type-I subtrees

Suppose that  $T = (\Lambda^{(2)}, \{\tau_1, \tau_2\}, \theta)$  is a Type-I subtree with the underlying model as described in Section 2. Let  $z_1$  and  $z_2$  be its two tips. Let the root be denoted as ‘0’ (Figure 1(d)) and let  $\tau_k$  be path distance between  $z_k$  and the root ( $k = 1, 2$ ).

**Proposition** Suppose that we have at least two lineages sampled at each of  $z_1$  and  $z_2$  and the root distribution is identifiable. Then  $\tau_1, \tau_2$  and  $\theta$  can be expressed as functions of the joint pmf of allele types in  $z_1$  and  $z_2$ , and hence they are identifiable.

**Proof** Suppose that we have samples of  $N_1$  and  $N_2$  lineages from  $z_1$  and  $z_2$  respectively, and the allele-counts among these lineages are  $R_1$  and  $R_2$  respectively. Let the joint pmf of  $(R_1, R_2)$  be  $f_{\mathbf{N}, \mathbf{R}}$ .

Consider random subsamples (without replacement) of size  $n_1$  and  $n_2$  from  $z_1$  and  $z_2$  respectively with  $n_k \leq 2, k = 1, 2$ . Rather than working with the allele-counts  $R_k$  at the original samples, we will work with allele-counts  $r_k$  at the subsamples.

One computes the joint pmf of  $(r_1, r_2)$  from  $f_{\mathbf{N}, \mathbf{R}}$  as

$$\begin{aligned} & \Pr_{\mathbf{n}}(r_k = j_k, k = 1, 2 \mid R_k = J_k, N_k = I_k, n_k = i_k, k = 1, 2; \tau_1, \tau_2, \theta) \\ &= \sum_{J_1=j_1}^{I_1-(i_1-j_1)} \sum_{J_2=j_2}^{I_2-(i_2-j_2)} \left( \prod_{k=1}^2 \frac{\binom{J_k}{j_k} \binom{I_k-J_k}{i_k-j_k}}{\binom{I_k}{i_k}} \right) f_{\mathbf{N}, \mathbf{R}}(J_1, J_2). \end{aligned}$$

We will argue that the joint pmfs  $(r_1, r_2)$  for  $(n_1, n_2) = (1, 1)$ ,  $(1, 2)$  and  $(2, 1)$  are enough to identify the parameters  $\tau_1, \tau_2$  and  $\theta$ .

As before, let  $n'_k$  be the number of lineages ancestral to subsamples at  $z_k$  that are present at the top node (the root) (see Figure 1(d)) and  $r'_k$  be the allele-count out of these  $n'_k$ ; ( $k = 1, 2$ ). Also, let  $n_0 = n'_1 + n'_2$  be the number of lineages at the root ancestral to the subsampled lineages at  $z_1$  and  $z_2$ , and  $r_0 = r'_1 + r'_2$  be the allele-count out of these  $n_0$  lineages.

First, consider the case  $n_1 = n_2 = 1$ . Then  $r_k = 0$  or  $1$  for  $k = 1, 2$ . From Eq. (2) it follows that  $n'_1 = n'_2 = 1$ ; thus,  $\Pr_{\mathbf{n}}(n'_k = i' \mid n_k = i; \tau_k)$  and hence

$\Pr(r_1 = j_1, r_2 = j_2 | n_1 = n_2 = 1; \tau_1, \tau_2, \boldsymbol{\theta})$  does not involve  $\tau_1$  and  $\tau_2$ . From Eq. (5) it also follows that  $r_k = r'_k, k = 1, 2$ . Also,  $n_0 = n'_1 + n'_2 = 1 + 1 = 2$ .

Note that  $r_0 = r'_1 + r'_2$  and  $r_k = r'_k (k = 1, 2)$  are counts. Thus,

$$(r_1, r_2) = (0, 0) \iff (r'_1, r'_2) = (0, 0) \iff r_0 = 0.$$

Using a symmetric argument

$$(r_1, r_2) = (1, 1) \iff (r'_1, r'_2) = (1, 1) \iff r_0 = 2.$$

Thus,

$$\Pr((r_1, r_2) = (j, j) | n_1 = n_2 = 1; \boldsymbol{\theta}) = \Pr(r_0 = 2j | n_0 = 2; \boldsymbol{\theta}), \quad j = 0, 1. \quad (7)$$

It follows that

$$\begin{aligned} & \Pr((r_1, r_2) = (0, 1) | n_1 = n_2 = 1; \boldsymbol{\theta}) + \Pr((r_1, r_2) = (1, 0) | n_1 = n_2 = 1; \boldsymbol{\theta}) \\ &= 1 - \Pr((r_1, r_2) = (0, 0) | n_1 = n_2 = 1; \boldsymbol{\theta}) - \Pr((r_1, r_2) = (1, 1) | n_1 = n_2 = 1; \boldsymbol{\theta}) \\ &= 1 - \Pr(r_0 = 0 | n_0 = 2; \boldsymbol{\theta}) - \Pr(r_0 = 2 | n_0 = 2; \boldsymbol{\theta}) \\ &= \Pr(r_0 = 1 | n_0 = 2; \boldsymbol{\theta}) \end{aligned} \quad (8)$$

Thus, from Eqs. (7) and (8)  $\Pr(r_0 = j_0 | n_0 = 2; \boldsymbol{\theta}), j_0 = 0, 1, 2$  can be expressed as functions of  $\Pr((r_1, r_2) = (j_1, j_2) | n_1 = n_2 = 1; \boldsymbol{\theta}), j_1, j_2 = 0, 1$ . The former is the root distribution for  $n_0 = 2$ , which is identifiable by the condition of Proposition 3.1. Thus,  $\boldsymbol{\theta}$  can be expressed as a function of the pmf of  $r_0$  (given  $n_0 = 2$ ), and thus as a function of joint pmf of  $(r_1, r_2)$ . Hence, it can also be expressed as a function of  $f_{\mathbf{N}, \mathbf{R}}$ .

Next, we consider  $n_1 = 2, n_2 = 1$ . Then  $r_1 = 0, 1$  or  $2$  and  $r_2 = 0$  or  $1$ . From Eq. (2) it follows that  $n'_2 = 1$ ; thus  $\Pr_{\mathbf{n}}(n'_2 = i'_2 | n_2 = i_2; \tau_2)$  and hence  $\Pr((r_1, r_2) = (0, 1) | n_1 = n_2 = 1; \tau_1, \tau_2, \boldsymbol{\theta})$  does not involve  $\tau_2$ . Moreover,  $n_0 = n'_1 + n'_2 = n'_1 + 1$ . Also, from Eq. (5) it follows that

$$(r_1, r_2) = (0, 1) \iff (r'_1, r'_2) = (0, 1).$$

Thus,

$$\begin{aligned} & \Pr((r_1, r_2) = (0, 1) | (n_1, n_2) = (2, 1); \tau_1, \boldsymbol{\theta}) \\ &= \sum_{i'=1}^2 \Pr((r'_1, r'_2) = (0, 1) | (n'_1, n'_2) = (i', 1); \boldsymbol{\theta}) \Pr(n'_1 = i' | n_1 = 2; \tau_1) \\ &= \sum_{i'=1}^2 \Pr((r'_1, r'_2) = (0, 1) | n_0 = n'_1 + 1 = i' + 1; \boldsymbol{\theta}) \Pr(n'_1 = i' | n_1 = 2; \tau_1) \\ &= \sum_{i'=1}^2 \sum_{j_0=0}^{i'+1} \Pr((r'_1, r'_2) = (0, 1) | r_0 = j_0, n_0 = i' + 1) \\ & \quad \times \Pr(r_0 = j_0 | n_0 = i' + 1; \boldsymbol{\theta}) \Pr(n'_1 = i' | n_1 = 2; \tau_1) \end{aligned}$$

Note that  $r_0 \neq 1 \implies (r'_1, r'_2) \neq (0, 1)$ . Also, note that

$$\Pr(r_0 = 1 \mid n_0 = i' + 1; \boldsymbol{\theta})$$

is a function of  $\boldsymbol{\theta}$  only (and no other parameters); hence we call it  $c_{i'+1}(\boldsymbol{\theta})$ ,  $i' = 1, 2$ . Thus,

$$\begin{aligned} & \Pr((r_1, r_2) = (0, 1) \mid (n_1, n_2) = (2, 1); \tau_1, \boldsymbol{\theta}) \\ &= \sum_{i'=1}^2 \Pr((r'_1, r'_2) = (0, 1) \mid r_0 = 1, n_0 = i' + 1) c_{i'+1}(\boldsymbol{\theta}) \Pr(n'_1 = i' \mid n_1 = 2; \tau_1) \\ &= \frac{c_2(\boldsymbol{\theta})}{2} (1 - e^{-\tau_1}) + \frac{c_3(\boldsymbol{\theta})}{3} e^{-\tau_1} = e^{-\tau_1} \left( \frac{c_3(\boldsymbol{\theta})}{3} - \frac{c_2(\boldsymbol{\theta})}{2} \right) + \frac{c_2(\boldsymbol{\theta})}{2} \end{aligned}$$

from Eqs. (2) and (4). From the above equation it follows that

$$\tau_1 = b\left(\Pr((r_1, r_2) = (0, 1) \mid (n_1, n_2) = (2, 1); \tau_1, \boldsymbol{\theta}), \boldsymbol{\theta}\right) \quad (9)$$

for some function  $b(\cdot, \cdot)$ . We have already established that  $\boldsymbol{\theta}$  can be expressed as a function of  $f_{\mathbf{N}, \mathbf{R}}$ . Thus,  $\tau_1$  can be expressed as a function of  $f_{\mathbf{N}, \mathbf{R}}$  and hence  $\tau_1$  is identifiable.

Using a symmetric argument, one can establish that  $\tau_2$  can be expressed as a function of  $f_{\mathbf{N}, \mathbf{R}}$  and hence it is identifiable. Thus, this proposition is proven.

### 3.2 Identifiability of Type-II subtrees

Consider a Type-II subtree of with tips  $z_A$  and  $z_B$ . Let the MRCA node of  $z_A$  and  $z_B$  be denoted as  $z_{AB}$ . (By definition  $z_{AB}$  is not the root.) Also, consider the path from  $z_{AB}$  to the root (node 0) and call it branch  $AB$ . There must be at least another branch  $H$  attached to the root other than branch  $AB$  (Figure 1(e)). Consider a tip  $z_D$ , such that the path between  $z_D$  and the root goes through  $H$ . Let  $\tau_A$  be the path distance between  $z_{AB}$  and  $z_A$  and let  $\tau_B$  be the path distance between  $z_{AB}$  and  $z_B$ . Also, let  $\tau_{AB}$  be the path distance between the root and  $z_{AB}$  and let  $\tau_D$  be the path distance between the root and  $z_D$ .

**Proposition** Suppose that we have at least two haploids sampled at each of  $z_A, z_B$  and  $z_D$  and the root distribution is identifiable. Then  $\tau_A, \tau_B, \tau_{AB}, \tau_D$  and  $\boldsymbol{\theta}$  can be expressed as functions of the joint pmf of the allele types at  $z_A, z_B$  and  $z_D$ , and hence they are identifiable.

**proof** Suppose that we have samples of  $N_A, N_B$  and  $N_D$  lineages from  $z_A, z_B$  and  $z_D$  respectively, and the allele-counts among these lineages are  $R_A, R_B$  and  $R_D$  respectively. Let the joint pmf of  $(R_A, R_B, R_D)$  be  $f_{\mathbf{N}, \mathbf{R}}^*$ .

First we consider the Type-I subtree formed by  $z_A$  and  $z_D$ . From Proposition 3.1 one can establish that  $\boldsymbol{\theta}, \tau_D$  and  $\tau_A + \tau_{AB}$  can be expressed as a function of the joint pmf of  $(R_A, R_D)$  and hence of  $f_{\mathbf{N}, \mathbf{R}}^*$ . A symmetric argument also establishes that  $\tau_B + \tau_{AB}$  can be expressed as functions of  $f_{\mathbf{N}, \mathbf{R}}^*$ . Next we will show that each of  $z_A, z_B$  and  $z_{AB}$  can be expressed as function of  $f_{\mathbf{N}, \mathbf{R}}^*$ .



Consider a random subsample of size one from each of  $z_A, z_B$  and  $z_D$ . Let  $n_A, n_B$  and  $n_D$  be the numbers of subsampled haploids at  $z_A, z_B$  and  $z_D$  respectively. (Thus,  $n_A = n_B = n_D = 1$ ). Let  $r_A, r_B$  and  $r_D$ , respectively, be the observed allele-counts at these subsamples. ( $r_k = 0$  or  $1$  for  $k = A, B, D$ .) As before, let  $n'_k$  be the number of lineages ancestral to subsamples at  $z_k$  that are present at the top node of the branch (in the subtree) attached to  $z_k$  (see Figure 1(e)) and  $r'_k$  be the allele-count out of these  $n'_k$  ( $k = A, B, D$ ).

From Eq. (2) it follows that  $n_A = n_B = n_D = n'_A = n'_B = n'_D = 1$  and thus  $\Pr(n'_k = i'_k | n_k = i_k; \tau_k)$  does not involve  $\tau_k$  ( $k = A, B, D$ ). Hence,

$$\Pr((r_A, r_B, r_D) = (0, 0, 1) | n_A = n_B = n_D = 1; \tau_A, \tau_B, \tau_{AB}, \tau_D, \boldsymbol{\theta})$$

does not involve  $\tau_A, \tau_B$  and  $\tau_D$ . Also,

$$\begin{aligned} & \Pr((r_A, r_B, r_D) = (0, 0, 1) | n_A = n_B = n_D = 1; \tau_{AB}, \boldsymbol{\theta}) \\ = & \sum_{J_A=j_A}^{I_A-(i_A-j_A)} \sum_{J_B=j_B}^{I_B-(i_B-j_B)} \sum_{J_C=j_C}^{I_C-(i_C-j_C)} \left( \prod_{k \in \{A, B, D\}} \frac{\binom{J_k}{j_k} \binom{I_k - J_k}{i_k - j_k}}{\binom{I_k}{i_k}} \right) f_{\mathbf{N}, \mathbf{R}}^*(J_A, J_B, J_D). \end{aligned} \quad (10)$$

Thus, the left side of Eq. (10) can be expressed as a function of  $f_{\mathbf{N}, \mathbf{R}}^*$ . It also follows from Eq. (5) that  $r_k = r'_k, k = A, B, D$ .

Let  $n_{AB} = n'_A + n'_B$  be the total number of lineages from subsamples of  $z_A$  and  $z_B$  that are present at node  $AB$ , and let  $r_{AB} = r'_A + r'_B$  be the allele-counts out of these  $n_{AB}$  lineages. Also, let  $n'_{AB}$  be the number of lineages ancestral to those  $n_{AB}$  lineages that are present at the top node (root) of the branch  $AB$ , and let  $r'_{AB}$  be the allele-count out of these  $n'_{AB}$  lineages. As before, let  $n_0 = n'_{AB} + n'_D$  be the total number of lineages at the root ancestral to the subsamples at  $z_A, z_B$  and  $z_D$ ; let  $r_0 = r'_{AB} + r'_D$  be the allele-count out of these  $n_0$  lineages. Note that  $n_{AB} = n'_A + n'_B = 2$ ,  $n'_{AB} \leq n_{AB}$ . From Eq. (5) and the fact that  $r_{AB} = r'_A + r'_B$  it follows that

$$(r_A, r_B) = (0, 0) \iff (r'_A, r'_B) = (0, 0) \iff r_{AB} = 0.$$

Thus,

$$\begin{aligned} & \Pr((r_A, r_B, r_D) = (0, 0, 1) | n_A = n_B = n_D = 1; \tau_{AB}, \boldsymbol{\theta}) \\ = & \Pr((r_{AB}, r_D) = (0, 1) | (n_{AB}, n_D) = (2, 1); \tau_{AB}, \boldsymbol{\theta}) \end{aligned} \quad (11)$$

Consider the part of the subtree consisting of the path from  $z_{AB}$  and  $z_D$  to the root; it is a Type-I subtree with  $z_{AB}$  and  $z_D$  as the tips, and  $\tau_{AB}$  and  $\tau_D$ , respectively, as the lengths of the attached branches; it has  $(n_{AB}, n_D) = (2, 0)$ , respectively, as the numbers of observed lineages at  $z_{AB}$  and  $z_D$  and  $(r_{AB}, r_D)$ , respectively, as the allele-counts in these lineages. From Eq. (9) and (11)

$$\begin{aligned} \tau_{AB} &= b \left( \Pr((r_{AB}, r_D) = (0, 1) | (n_{AB}, n_D) = (2, 1); \tau_{AB}, \boldsymbol{\theta}), \boldsymbol{\theta} \right) \\ &= b \left( \Pr((r_A, r_B, r_D) = (0, 0, 1) | n_A = n_B = n_D = 1; \tau_{AB}, \boldsymbol{\theta}), \boldsymbol{\theta} \right). \end{aligned}$$

As we have already established that  $\tau_A + \tau_{AB}$ ,  $\tau_B + \tau_{AB}$ ,  $\tau_D$ ,  $\theta$  and the left side of Eq. (10) can be expressed as functions of  $f_{\mathbf{N},\mathbf{R}}^*$ , it follows that  $\tau_A, \tau_B, \tau_{AB}, \tau_D$  and  $\theta$  can be expressed as functions of  $f_{\mathbf{N},\mathbf{R}}^*$ . Thus, they are identifiable and this proposition is proven.

Thus, the parameters of the tree are identifiable, as each two-tip subtree along with the root distribution parameter  $\theta$  is identifiable.

## 4 Discussions

We have proven that the model parameters are identifiable under the coalescent-based population tree model of [2, 5]. Thus, the problem of estimation of population tree from this model is indeed meaningfully stated. Moreover, as identifiability is a required condition for consistency of maximum likelihood estimator (MLE), this is a step towards proving the consistency of MLE for this model. We have proven the identifiability of the tree parameters for any identifiable root distribution. As a result our proof is valid for different versions of this model (that vary at the root distribution) such as [2, 5, 6].

## References

- [1] ALLMAN, E. S., ANÉ, C. & RHODES, J. A. (2008). Identifiability of a Markovian model of molecular evolution with gamma-distributed rates. *Adv. Appl. Prob.* **40**, 228–249.
- [2] NIELSEN, R., MOUNTAIN, J. L., HUELSENBECK, J. P. & SLATKIN, M. (1998). Maximum likelihood estimation of population divergence times and population phylogeny in models without mutation. *Evolution* **52**, 669–677.
- [3] J. F. C. KINGMAN. The coalescent. (1982) *Stoch. Proc. Applns.*, 13:235–248, 1982.
- [4] R. Nielsen & M. Slatkin. (2000) Likelihood analysis of ongoing gene flow and historical association. *Evolution*, 54:44–50.
- [5] ROYCHOUDHURY, A., FELSENSTEIN, J. & THOMPSON, E. A. (2008). A two-stage pruning algorithm for likelihood computation for a population tree. *Genetics* **180**. 1095–1105.
- [6] ROYCHOUDHURY, A. (2011). Composite likelihood-based inferences on genetic data from dependent loci. *Journal of Mathematical Biology* **180**. 62(1):65–80.
- [7] ROYCHOUDHURY, A., & THOMPSON, E. A. (2012). Ascertainment correction for a population tree via a pruning algorithm for likelihood computation. *Theoretical Population Biology* **180**. 82(1):59–65.
- [8] BRYANT, D., BOUCKAERT, R., FELSENSTEIN, J., ROSENBERG, N. A. & ROYCHOUDHURY, A. (2012). Inferring species trees directly from biallelic genetic markers: bypassing gene trees in a full coalescent analysis. *Mol. Biol. Evol.* **29**, 1917–1932.
- [9] CHAI, J. & HOUSWORTH, E. A. (2010). On Rogers’ proof of identifiability for the GTR +  $\Gamma$  + I Model. *Syst. Biol* **60**, 713–718.
- [10] STEEL, M. A., SZÉKERLY, L. A. & HENDY, M. D. (1994). Reconstructing trees when sequence sites evolve at variable rates. *J. Comp. Biology* 1(2):153–163.

- [11] MATSON, F. A., & STEEL, M. (2008). Phylogenetic mixtures on a single tree can mimic a tree of another topology. *Syst. Biology* 56(5):767-775.
- [12] FELSENSTEIN, J. (1981). Evolutionary trees from DNA sequences: a maximum likelihood approach. *J. Mol. Evol.* **17**, 368–376.
- [13] L. Liu, L. Yu, D. K. Pearl & S. V. Edwards (2009). Estimating species phylogenies using coalescence times among sequences. *Systematic Biology*, 58(5):468–477.